Systematic Probing for Examining Evidence Based Tabular Reasoning

https://tabevidence.github.io/

Vivek Gupta^{1*}, Riyaz Bhat², Atreya Ghosal³, Manish Shrivastava³ Maneesh Singh², Vivek Srikumar¹

¹University of Utah; Verisk Analytics²; IIIT Hyderabad³ INTERNATIONAL INSTITUTE O *Bloomberg Ph.D. Fellow (2021-2022)

Bloomberg

*on academic job market



Any *"evidence-based reasoning"* system should demonstrate *expected*, *predictable behavior* in *response* to *controlled changes to its inputs*.

Case Study on Tabular Inference

TAKEAWAY

Systematic target probing can highlights the limitation of tabular reasoning models

Such targeted probes are data efficient and work with minimal to no supervision

Probing can highlights several problems in tabular reasoning models:

(a) Use of annotation artifacts

(b) Use of incorrect evidence

(c) Overfitting on pre-trained knowledge

Check out the paper for details: <u>https://tabprobe.github.io</u>

TABULAR INFERENCE

- The **tabular natural language inference** problem is similar to standard NLI
- But here, the **premises are tabular data**
- Task: to decide whether given hypothesis is true (entailment), false (contradiction) or undetermined (neutral) given a premise table

Check out InfoTabS (Gupta et al., 2020) https://infotabs.github.io

Drobbugo		
Highest governing body	International Federation for Equestrian Sports (FEI)	
	Characteristics	
Contact	No	
Team members	Individual and team at international levels	
Mixed gender	Yes	
Equipment	Horse, appropriate horse tack	
Venue	Arena, indoor or outdoor	
Presence		
Country or region	Worldwide	
Olympic	1912	
Paralympic	1996	

Dreesana

H1: Both men and women can complete in the contactless sport of Dressage \rightarrow Entail

EVIDENCE BASED REASONING

But why entail? assume that model is

 \rightarrow attention at **evidence** (relevant rows) \rightarrow correct **reasoning** (logical inference)

<u>black box problem</u>

"model doesn't provide the highlighted evidence and the reasoning steps"

Highest governing body	International Federation for Equestrian Sports (FEI)	
	Characteristics	
Contact	No	
Team	Individual and team at	
members	international levels	
Mixed gender	Yes	
Equipment	Horse, appropriate horse tack	
Venue Arena, indoor or outdoor		
	Presence	
Country or region	Worldwide	
Olympic	1912	
Paralympic	1996	

Dressage **CONTROLLED** CHANGES Highest International Federation for governing Both men and women can compete in the Equestrian Sports (FEI) body contactless sport of Dressage **Characteristics** Contact No Individual and team at Team remove the row with the key international levels members contact -Mixed gender Yes Equipment Horse, appropriate horse tack model prediction change Arena, indoor or outdoor Venue entail \rightarrow ? Presence **Country or** Worldwide region SYSTEMATIC PROBES Olympic 1912 1996 Paralympic

EXPECTED MODEL RESPONSE

<u>EXPECTED</u> WODEL RESPO	NSE	Dressage	
remove the key contact from the table	Highest governing body	International Federation for Equestrian Sports (FEI)	
Both men and women can compete in the	Characteristics		
contactless sport of Dressage	Team members	Individual and team at international levels	
deally prediction change	Mixed gender	Yes	
entail \rightarrow neutral	Equipment	Horse, appropriate horse tack	
	Venue	Arena, indoor or outdoor	
(c	Presence		
	Country or region	Worldwide	
	Olympic	1912	
	Paralympic	1996	

Systematic ('CONTROLLED') **Probes** (CHANGES)

In our study, we define **three types** of systematic probes, as follow:

1. Annotation Artifacts:

Can a model make inference about a hypothesis without a premise (a.k.a using artifacts)?

ANNOTATION ARTIFACTS

Can a model make inference about a hypothesis without a premise (a.k.a using artifacts)?

Modify the hypothesis in such a way that the inference label is retained or flipped

ANNOTATION ARTIFACTS

Modify the following types of expressions to alter the hypothesis:

Named Entity	involving named entities		
Numerical	related to numbers	for others refer to the peper	
Temporal	involving date and time	for others refer to the paper	
Quantification	related to introduction, deletion or modification of quantifiers such as most, many, every, etc.		
Lexical	lexical semantics i.e. antonymy, synonymy, etc.		
Negation	introduction or deletion of negation markers		
Syntactic Alternations	leveraging syntactic structure		
Subjective	adding subjective phrases/expressions		

EXAMPLE: ANNOTATION ARTIFACTS

Named Entity

H: Katie Holmes moved from **Ohio To California**.

Prediction: Entail



H': Katie Holmes moved from **South Africa** To California.

Prediction: Contradiction

_	
ti	tle
к	Katie Holmes
В	orn
(1978-12-18) December 18, 1978 (age 40) Toledo, Ohio, U.S.
В	irth name
K	Cate Noelle Holmes
R	esidence
L	os Angeles, California, U.S.
0	ther names
K	Katherine Noelle Holmes
0	occupation
A	Actress
Y	ears active
1	.997-present
S	pouse(s)
Т	om Cruise (m. 2006 ; div. 2012)
С	hildren
1	

	×	
$\begin{array}{c} \textbf{Original} \\ mean_{(stdev)} \end{array}$	Label Pre- served	Label Flipped
Train Set (w/o	NEUTRAL)	
99.44 (0.06)	92.98 _(0.20)	53.92(0.28)
96.39 _(0.13)	$70.23_{(0.35)}$	19.23(0.27)
α_1 Set (w/o N	NEUTRAL)	
68.94 _(0.76)	69.56 (0.77)	51.48(0.86)
63.52(0.75)	60.27(0.85)	31.02(0.63)
	Original $mean_{(stdev)}$ $Train$ Set (w/o 99.44 _(0.06) 96.39 _(0.13) α_1 Set (w/o N 68.94 _(0.76) 63.52 _(0.75)	Original mean _(stdev) Label Pre- served Train Set (w/o NEUTRAL) 99.44 _(0.06) 92.98 _(0.20) 96.39 _(0.13) 70.23 _(0.35) α ₁ Set (w/o NEUTRAL) 68.94 _(0.76) 69.56 _(0.77) 63.52 _(0.75) 60.27 _(0.85)

Model retain performance when
label is preserve after perturbation

We tried two kinds of perturbations setting:

- a) that preserve the label,
- b) and those that flipped the label from entail to contradict and vice versa

Model	$\begin{array}{c} \textbf{Original} \\ mean_{(stdev)} \end{array}$	Label Pre- served	Label Flipped
5	Train Set (w/o	NEUTRAL)	
Prem+Hypo	99.44 _(0.06)	92.98 (0.20)	53.92 _(0.28)
Hypo-Only	96.39 _(0.13)	70.23(0.35)	19.23(0.27)
	α_1 Set (w/o N	NEUTRAL)	
Prem+Hypo	$68.94_{(0.76)}$	69.56 _(0.77)	$51.48_{(0.86)}$
Hypo-Only	$63.52_{(0.75)}$	$60.27_{(0.85)}$	$31.02_{(0.63)}$

1. Model retain performance when label is preserve after perturbation

 Substantial drop only for Hypothesis only baseline on training settings.

Model	$\begin{array}{c} \textbf{Original} \\ mean_{(stdev)} \end{array}$	Label Pre- served	Label Flipped
9 11	Train Set (w/o	NEUTRAL)	
Prem+Hypo	99.44 (0.06)	92.98 _(0.20)	$53.92_{(0.28)}$
Hypo-Only	96.39 _(0.13)	70.23(0.35)	15.23(0.27)
α_1 Set (w/o NEUTRAL)			
Prem+Hypo	$68.94_{(0.76)}$	69.56 _(0.77)	$51.48_{(0.86)}$
Hypo-Only	$63.52_{(0.75)}$	$60.27_{(0.85)}$	$31.02_{(0.03)}$

1. Model retain performance when label is preserve after perturbation

- 2. Substantial drop only for Hypothesis only baseline on training settings.
- 3. Model overfit on hypothesis baseline in original setting

Model	Original	Label Pre-	Label
	$mean_{(stdev)}$	served	Flipped
5 15	Train Set (w/o	NEUTRAL)	
Prem+Hypo	99.44 _(0.06)	92.98 _(0.20)	53.92 _(0.28)
Hypo-Only	96.39 _(0.13)	$70.23_{(0.35)}$	$19.23_{(0.27)}$
	α_1 Set (w/o I	NEUTRAL)	
Prem+Hypo	68.94 _(0.76)	69.56 _(0.77)	51.48(0.86)
Hypo-Only	$63.52_{(0.75)}$	$60.27_{(0.85)}$	31.02(0.63)
	k		
			<u>`</u>

- 1. Model retain performance when label is preserve after perturbation
- 2. Substantial drop only for Hypothesis only baseline on training settings.
- 3. Model overfit on hypothesis baseline in original setting.
- 4. Model performance dropped drastically when label is flipped after perturbation

FINDINGS: ANNOTATION ARTIFACTS

Can a model make inference about a hypothesis without a premise (a.k.a using artifacts)?

Yes, models largely rely on spurious correlation between hypothesis sentence and inference label.

for detailed results check the paper

Evidence Selection:

In our study, we define **three types** of systematic probes, as follow:

- 1. Annotation Artifacts:
 - a. Can a model make inference about a hypothesis without a premise (a.k.a using artifacts)?

2. Evidence Selection:

a. Is the model drawing inferences based on right evidence in the premise?

EVIDENCE SELECTION

Is the model drawing inferences based on right evidence in the premise?

Systematically alter the premise table via simple operations in order to deterministically change the inference label.

 \rightarrow invalid label changes can be deterministically identified

POSSIBLE OPERATIONS

- Row Deletion
 - Automatic Probing: Random Row Deletion
 - Manual Probing: Require Row Annotation
 - Relevant Row Deletion
 - Irrelevant Row Deletion
- Row-Value Update
- New Row Insertion
- Row Perturbation

ANY ROW DELETION



RANDOM ROW DELETION



**RoBERTa Large from Gupta, Vivek, et al. "INFOTABS: Inference on Tables as Semi-structured Data." ACL 2020.

RELEVANT ROW DELETION

remove the key "contact" from the table	Highest governing body	International Federation for Equestrian Sports (FEI)
Both men and women can play in the contactless sport	Cł	naracteristics
of Dressage.	Team members	Individual and team at international levels
Prediction: Neutral	Mixed gender	Yes
	Equipment	Horse, appropriate horse tack
	Venue	Arena, indoor or outdoor
		Presence
	Country or region	Worldwide
	Olympic	1912
	Paralympic	1996

Dressage

RELEVANT ROW DELETION

			Diessage
remove the key "contact"	from the table	Highest governing body	International Federation for Equestrian Sports (FEI)
Both men and women car	elay in the contactless sport	CI	naracteristics
of Dressage.		Team members	Individual and team at international levels
Prediction: Neutral	``>	Mixed gender	Yes
		Equipment	Horse, appropriate horse tack
E	N C	Venue	Arena, indoor or outdoor
E(381)	188 174		Presence
N 7	625 37	Country or region	Worldwide
C 71	175 680	Olympic	1912
		Paralympic	1996

Dragage

ROW INSERTION

Insert row "Arena size" in the table

Both Men and Women can play in the contactless sport of Dressage.

Prediction : ??

	Dressage
Highest governing body	International Federation for Equestrian Sports (FEI)
CI	haracteristics
Contact	No
Team members	Individual and team at international levels
Mixed gender	Yes
Equipment	Horse, appropriate horse tack
Venue	Arena, indoor or outdoor
Arena Size	20 by 60 m [66 by 197 ft]
	Presence
Country or region	Worldwide
Olympic	1912
Paralympic	1996

			Dressage		
	NOLN			Highest governing body	International Federation for Equestrian Sports (FEI)
insert row	Arena size	"in the table		CI	naracteristics
			L>	Contact	No
Both Men a of Dressage	and Wome e.	en can p lay in	the contactless sport	Team members	Individual and team at international levels
			· · · >	Mixed gender	Yes
Prediction :	Entail	C	\bigcirc	Equipment	Horse, appropriate horse tack
	106	367		Venue	Arena, indoor or outdoor
N 411	14168	294	C	Arena Size	20 by 60 m [66 by 197 ft]
C <mark>695</mark>	355	14920		Presence	
		\bigcirc		Country or region	Worldwide
		V е		Olympic	1912
				Paralympic	1996
			\sim		1.1 (1997) 1997 - State State (1997) 1997 - State State (1997)

FINDING: EVIDENCE SELECTION

Is the model drawing inferences based on right evidence in the premise?

for detailed results check the paper

Systematic ('CONTROL') **Probes** (CHANGES):

In our study, we define **three types** of systematic probes, as follow:

1. Annotation Artifacts:

Can a model make inference about a hypothesis without a premise (a.k.a using artifacts)?

2. Evidence Selection:

Is the model drawing inferences based on right evidence in the premise?

3. Counterfactual Instances:

How will the model react if the primary evidence is counterfactual to pre-trained data?

COUNTERFACTUAL INSTANCES

How will the model react if the primary evidence is counterfactual to pre-trained data?

Update the premise table to include counterfactual data in order to retain or change the inference label.

EXAMPLE: COUNTERFACTUAL INSTANCES

Dressage				
Highest governing body	International Federation for Equestrian Sports (FEI)			
Characteristics				
Contact	No			
Team members	Individual and team at international levels			
Mixed gender	Yes			
Equipment	Horse, appropriate horse tack			
Venue	Arena, indoor or outdoor			
Presence				
Country or region	Worldwide			
Olympic	1912			
Paralympic	1996			

Both men and women can complete in the **contactless** sport of Dressage \rightarrow Entail

Dressage				
Highest governing body	International Federation for Equestrian Sports (FEI)			
Characteristics				
Contact	Yes			
Team members	Individual and team at international levels			
Mixed gender	Yes			
Equipment	Horse, appropriate horse tack			
Venue	Arena, indoor or outdoor			
Presence				
Country or region	Worldwide			
Olympic	1912			
Paralympic	1996			

Both men and women can complete in the **contactless** sport of Dressage \rightarrow **Eptatiladict**

EXAMPLE: COUNTERFACTUAL INSTANCES

Dressage				
Highest governing body	International Federation for Equestrian Sports (FEI)			
Characteristics				
Contact	No			
Team members	Individual and team at international levels			
Mixed gender	Yes			
Equipment	Horse, appropriate horse tack			
Venue	Arena, indoor or outdoor			
Presence				
Country or region	Worldwide			
Olympic	1912			
Paralympic	1996			

Both men and women can complete in the contactless sport of Dressage \rightarrow Entail

Dressage			
Highest governing body	International Federation for Equestrian Sports (FEI)		
Characteristics			
Contact	No		
Team members	Individual and team at international levels		
Mixed gender	No		
Equipment	Horse, appropriate horse tack		
Venue	Arena, indoor or outdoor		
Presence			
Country or region	Worldwide		
Olympic	1912		
Paralympic	1996		

Both men and women can complete in the contactless sport of Dressage \rightarrow **Eptatiladict**

FINDINGS: ANNOTATION ARTIFACTS

How will the model react if the primary evidence is counterfactual to pre-trained data?

Model relies on information from pre-trained language models rather than tabular evidence for making prediction

for detailed results check the paper

MAIN FINDINGS

- Annotation Artifacts: Model largely rely on spurious correlation between hypothesis and inference label.
- Evidence Selection: The model does not look at correct evidence required for correct reasoning.
- **Counterfactual Instances:** Model relies on information from pre-trained language model rather than tabular evidence.
- **Inoculation Study:** Changes in the data distribution during training have a negative impact on model performance.

TAKEAWAY

Systematic target probing can highlight the limitation of tabular reasoning models

Such targeted probes are data efficient and work with minimal to no supervision

Probing can highlights problems in tabular reasoning models:

(a) Use of annotation artifacts

(b) Use of incorrect evidence

(c) Overfitting on pre-trained knowledge

Check out the paper for details: <u>https://tabprobe.github.io</u>